

PRIORITY DOCUMENT
SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH
RULE 17.1(a) OR (b)



REC'D 10 SEP 2003

WIPO PCT

**Prioritätsbescheinigung über die Einreichung
einer Patentanmeldung**

Aktenzeichen: 102 29 207.8

Anmeldetag: 28. Juni 2002

Anmelder/Inhaber: T-Mobile Deutschland GmbH, Bonn/DE

Bezeichnung: Verfahren zur natürlichen Spracherkennung
auf Basis einer Generativen Transformations-
/Phrasenstruktur-Grammatik

IPC: G 10 L 15/18

**Die angehefteten Stücke sind eine richtige und genaue Wiedergabe der ur-
sprünglichen Unterlagen dieser Patentanmeldung.**

München, den 28. August 2003
Deutsches Patent- und Markenamt
Der Präsident
Im Auftrag

Stemme

28.06.2002

T-Mobile Deutschland GmbH

Verfahren zur natürlichen Spracherkennung auf Basis einer Generativen Transformations-/Phrasenstruktur-Grammatik

Die Erfindung betrifft ein Verfahren zur natürlichen Spracherkennung auf Basis einer Generativen Transformations-/Phrasenstruktur-Grammatik (GT/PS-Grammar).

Aktuelle Spracherkennungssysteme mit natürlicher Spracherkennung (NLU = Natural Language Understanding) sind in der Lage, eine Vielzahl möglicher Äußerungen zu verstehen und in komplexe Befehlsstrukturen umzusetzen, die Spracherkennungssysteme, z.B. Computer, zu bestimmten Aktionen veranlassen. Sie tun dies auf der Grundlage vorab definierter sinnvoller Mustersätze, die von Applikationsentwicklern und sogenannten Dialog-Designern festgelegt werden. Diese Sammlung von Mustersätzen – auch „Grammar“ genannt – umfasst einzelne Kommandoworte ebenso wie komplizierte Schachtelsätze, die an einer bestimmten Stelle des Dialogs sinnvoll sind. Äußert der Nutzer einen solchen Satz, wird er vom System mit großer Sicherheit verstanden und die mit ihm verknüpfte Handlungsanweisung wird ausgeführt.

Bei der Programmierung einer Erkennungsapplikation, z.B. einer NLU-Telefonapplikation, ist die Grammar also ein unverzichtbarer Baustein. Sie wird mit Hilfe eines speziellen Werkzeugs, der sogenannten Grammar Specification Language (GSL) erzeugt. Mit ihr werden die zu verstehenden Worte sowie ihre Verknüpfungen vorab reproduziert und für den Spracherkenner festgeschrieben. Die vorgegebenen Sätze werden dabei aus Wortkombinationen gebildet, die untereinander austauschbar (paradigmatische Achse) und miteinander kombinierbar (syntagmatische Achse) sind. Ein Beispiel hierfür ist in Figur 7 dargestellt.

Die möglichen Äußerungen ergeben sich aus der syntagmatischen Verknüpfung der paradigmatischen Wortkombinationen. Dass dabei auch Sätze möglich werden, die grammatisch falsch sind, wie z.B. „Würden Sie vielleicht Telly-Tarif ersetzen?“ muss in Kauf genommen werden, um das Antwortenspektrum möglichst groß zu halten. Diese sogenannte „Overgeneration“, das heißt z.B. das Vorhalten bzw. Erkennen von unsinnigen Mustersätzen oder Ausdrücken mit dem selben Sinngehalt, sollte jedoch gering gehalten werden, denn sie beansprucht beträchtliche Systemressourcen und setzt gleichzeitig die Erkennungsleistung herab, weil das System jede Nutzeräußerung mit einer Fülle vorgegebener Satzkombination vergleichen muss, die kaum jemals geäußert werden.

In der bisher üblichen Praxis wurden die paradigmatischen Wortkombinationen in einer Weise festgelegt, die scheinbar Zusammengehöriges verbindet. Dabei wurde von der bedeutungstragenden Qualität der Worte ausgegangen. Dieses Verfahren, das von einem mutmaßlichen Erfolgssatz ausgeht, entspricht durchaus den Erfordernissen einfacher Applikationen und führt hier zu zufriedenstellenden Ergebnissen. Bei komplexen Anwendungen, mit einer Fülle sinnvoller Antwortmöglichkeiten hingegen, werden diese herkömmlichen Grammatiken so groß, dass sie selbst die Rechenkapazität gegenwärtiger Hochleistungsserver bis an die Grenze belasten. Die Folgen sind:

- Stark vermehrte Overgeneration
- Spürbare Verzögerungen bei der Erkennung (Latency)
- Sinkende Erkennungssicherheit (Accuracy).
- Abgesenkte Systemstabilität (Robustness)

Der Hauptmangel dieser Methode besteht darin, dass die spezifizierten Sätze lediglich einer oberflächlichen Kombinatorik folgen. Die erzeugte Overgeneration ist deshalb so groß, weil die scheinbar zusammengehörigen Elemente tatsächlich anderen Kombinationsregeln folgen, die in der Sprachwissenschaft seit längerem bekannt sind.

Zusammenfassend wird festgehalten, dass die derzeit verbreiteten Grammars, die festlegen, welche Sätze von einem ASR-System erkannt werden, traditionellen grammatischen Konventionen folgen, die natürlich-sprachliche Äußerungen unzureichend strukturiert abbilden. Dabei wurde bislang nicht von einer Differenzierung von „Oberflächen-“, bzw. „Tiefenstrukturen“ ausgegangen. Die linguistische Hypothese besagt, dass eine syntaktische Tiefenstruktur und deren „generative Umsetzung“ hin zu konkreten Oberflächenstrukturen die Leistungsfähigkeit eines Sprachsystems ausmacht. Wird bei steigender Komplexität ausschließlich die bisher eingesetzte Oberflächenstruktur verwendet, muss diese, um ihrer Aufgabe dennoch gerecht zu werden, so groß dimensioniert sein, dass sie im Betrieb kaum noch vernünftig gepflegt werden kann und die Server bis an die Grenzen ihrer Kapazität belastet.

Die Aufgabe der Erfindung besteht darin, ein Verfahren zur Spracherkennung auf Basis einer Generativen Transformations-/Phrasenstruktur-Grammatik anzugeben, das im Vergleich zu herkömmlichen Erkennungsverfahren weniger Systemressourcen benötigt und dadurch eine sichere und schnelle Erkennung von Sprache bei gleichzeitiger Verringerung der Overgeneration ermöglicht.

Diese Aufgabe wird erfindungsgemäß durch die Merkmale des Patentanspruchs 1 gelöst.

Erfindungsgemäß erfolgt eine Analyse einer gesprochenen Phrase auf darin enthaltene Triphone, eine Bildung von in der gesprochenen Phrase enthaltenen Wörtern aus den erkannten Triphonen mit Hilfe von Lautwortdatenbasen (Dictionaries) und eine syntaktische Rekonstruktion der gesprochenen Phrase aus den erkannten Wörtern unter Verwendung eines grammatischen Regelwerks (Grammar).

Vorteilhafte Ausgestaltungen und Weiterbildungen der Erfindung ergeben sich aus den Merkmalen der Unteransprüche.

Besonders markant ist der Gegensatz zwischen dem erfindungsgemäßen Verfahren und der traditionellen Grammar Specification Language, die bei kleinen Applikationen auch mit syntaktischen Oberflächen, d.h. konkretes Ausformulieren von Erfolgssätzen, gute Resultate erzielte.

Erfindungsgemäß werden die Verknüpfungsregeln grammatischer Sätze nicht an der Oberfläche reproduziert, sondern die Tiefenstrukturen aufgezeigt, denen die syntagmatischen Verknüpfungen aller indogermanischen Sprachen folgen. Jeder Satz wird anhand eines syntaktischen Modells in Form von sogenannten Strukturbäumen beschrieben.

Die GT/PS-Grammar orientiert sich nicht an den potenziellen Äußerungen einer spezifischen Applikation, sondern an der Tiefenstruktur der Syntax (Satzbildungsregeln) indogermanischer Sprachen. Sie liefert ein Gerüst, das mit verschiedenen Worten gefüllt werden kann und die Realität der gesprochenen Sprache besser abbildet, als das bisher praktizierte „mimetische“ Verfahren.

Innerhalb der durch die Strukturbäume beschriebenen Tiefenstrukturen wird erkennbar, dass sich bestimmte Phrasen innerhalb eines Satzes wiederholen. Solche Wiederholungen können mit Hilfe der GSL reproduziert und aufgefangen werden. Dadurch sinkt nicht nur der Umfang einer Grammar erheblich, sondern auch die Overgeneration von grammatisch inkorrekten Sätzen sinkt beträchtlich.

Während in der traditionellen GSL-Grammar z.B. rund 500 Subgrammars in sieben hierarchischen Ebenen miteinander verflochten sind, kann die Anzahl der Subgrammars im GT/PS-Modell auf z.B. 30 Subgrammars in nur zwei hierarchischen Ebenen reduziert werden.

Der neue Grammartyp bildet natürlich-sprachliche Äußerungen in strukturierter Form ab und hat dabei z.B. nur rund 25% der Größe der bisherigen Grammar. Aufgrund ihrer geringen Größe ist diese Grammar einfacher zu pflegen, wobei die Zeiten für Kompilierung rapide sinken. Aufgrund ihrer geringen Größe steigt die

Erkennungssicherheit (Accuracy) und sinkt die Erkennungsverzögerung (Latency). Die aktuellen Rechnerkapazitäten werden besser ausgenutzt und die Performance der Server steigt. Darüber hinaus ist die neue Grammar nicht auf eine bestimmte Applikation bezogen, sondern kann in ihren Grundstrukturen für unterschiedliche Anwendungen verwendet werden, wodurch die Homogenität der Systeme gesteigert und die Entwicklungszeiten reduziert werden.

Der universale Code der Tiefenstruktur ermöglicht den Einsatz und die Wertschöpfung für multilinguale Sprachsysteme in einer bislang nicht erreichten Dimension, besonders die westeuropäischen Standardsprachen können mit vergleichsweise geringem Aufwand verarbeitet werden.

Im Unterschied zur bisherigen Grammar für natürlich-sprachliche Dialogapplikationen basiert die neue GT/PS Grammar auf aktuellen sprachwissenschaftlichen Modellen, die natürlich-sprachliche Äußerungen im Rahmen von Oberflächen- und Tiefenstrukturen abbilden. Die abstrakten Strukturmuster werden mit einer Grammar Specification Language (GSL) in ein hierarchisch verschachteltes und vernetztes Regelwerk übertragen, dessen Strukturen in der beiden Anlagen abgebildet sind.

Die technischen Vorzüge der GT/PS-Grammar sind damit:

- Die GT/PS-Grammar ist sehr viel kleiner als die bisherige Grammar, weil sie statt der bisher bis zu sieben Subgrammarlevels nur noch mit zwei Ebenen auskommt;
- Die Zahl der von der Grammar abgedeckten aber grammatisch falschen Sätze (Overgeneration) sinkt drastisch;
- Sie benötigt nur noch rund ein Drittel der bislang verwendeten Slots;
- Sie füllt entgegen der heutigen Spracherkenner-Philosophie die Slots in den unteren Grammar-Ebenen, statt in den oberen;
- Sie nutzt das von der GSL (Grammar Specification Language) bereit gestellte Instrument, Slotwerte in höhere Grammarlevels hoch zu reichen, konsequent aus;

- Sie besitzt einen neuen Slot mit der Bezeichnung ACTION, der nur noch mit den Werten GET und KILL gefüllt werden kann;
- sie arbeitet mit ineinander verschachtelten Slots, die hochgradig multitaskingfähig sind.
- Sie führt zu einer Verbesserung der Erkennerleistung
- Sie ermöglicht eine vereinfachte Option zur Einführung mehrsprachiger Applikationen
- Sie weist eine nahtlose Integrationsfähigkeit in Nuance Technologie auf

Die wirtschaftlichen Vorzüge der PSG sind:

- Verringerung der Hardwarekosten durch bessere Ausnutzung der Systemressourcen
- Verringerung der Übertragungszeiten durch leistungsfähigere Erkennung
- Einsparung von Personalressourcen durch leichtere Pflegbarkeit
- Größere Kundenzufriedenheit
- Anwendbar auf alle Weltsprachen (Englisch bis Chinesisch)

Nachfolgend wird die Erfindung anhand eines vereinfachten Ausführungsbeispiels unter Bezugnahme auf die Zeichnungen näher erläutert. Aus den Zeichnungen und deren Beschreibung ergeben sich weitere Merkmale, Vorteile und Anwendungsmöglichkeiten der Erfindung.

Es zeigt.

- Figur 1: Eine Triphonanalyse als ersten Schritt im Erkennungsprozess;
- Figur 2: Eine Worterkennung aus den erkannten Triphonen als zweiten Schritt im Erkennungsprozess;
- Figur 3;: eine syntaktische Rekonstruktion der erkannten Wörter als dritten Schritt des Erkennungsprozesses;
- Figur 4: Ein Beispiel für die Gliederung der erkannten Wörter in Wortartenkategorien sowie in nominale und verbale Phrasen;
- Figur 5: Ein Programmbeispiel für eine mögliche Grammar;

Figur 6: Eine Übersicht über die Struktur einer PSG Grammar;

Figur 7: Ein Beispiel für eine Bildung von Wortkombinationen bei einer Grammar nach den Stand der Technik.

Figur 1 zeigt den ersten Schritt einer Spracherkennung: die Triphonanalyse. Der kontinuierliche Redefluss einer Person 1 wird z.B. von einem Mikrofon eines Telefons angenommen und als analoges Signal einem Spracherkenner 2 zugeführt. Dort wird das analoge Sprachsignal in ein digitales Sprachsignal 3 umgewandelt. Das Sprachsignal enthält eine Vielzahl von Triphonen, d.h. Lautsegmenten, die im Spracherkenner 2 mit vorhandenen, d.h. vorgegebenen Triphon-Verknüpfungsregeln abgeglichen werden. Die vorhandenen Triphone sind in einer Datenbasis abgespeichert, die ein oder mehrere Lautwörterbücher enthält. Die erkannten Triphone liegen dann als eine Triphon-Kette 4 vor, z.B. „pro“, „rot“, „ote“, „tel“.

In einem zweiten Schritt gemäß Figur 2 werden aus den erkannten Triphonen sinnvolle Wörter gebildet. Dazu wird die vorhandene Triphon-Kette 4 mit in einem weiteren Lautwörterbuch 5 abgespeicherten, vorgegebenen Wörtern 6, z.B. „profi“, „portal“, „protel“, „hotel“, verglichen. Das Lautwörterbuch 5 kann einen bestimmten Wortschatz aus der Umgangssprache sowie einen auf die jeweilige Anwendung zugeschnittenen, speziellen Wortschatz umfassen. Stimmen die erkannten Triphone, z.B. „pro“ und „tel“, mit den in einem Wort, z.B. „protel“, enthaltenen Triphonen überein, wird das entsprechende Wort 7 als solches erkannt: „protel“.

Im nächsten Schritt, dargestellt in Figur 3, erfolgt die syntaktische Rekonstruktion der erkannten Wörter 7 mit Hilfe der Grammar 8. Dazu werden die erkannten Wörter ihren Wortartkategorien, wie Nomen, Verb, Adverb, Artikel, Adjektiv, etc. zugeordnet, wie dies in Figur 6 dargestellt ist. Dies erfolgt anhand von in Wortartkategorien unterteilten Datenbasen. Wie man in Figur 5 erkennt, können die Datenbasen 9-15 sowohl die oben erwähnten, herkömmlichen Wortartkategorien als auch spezielle Wortartkategorien, wie z.B. Ja/Nein Grammtik 9, Telefonnummern 14, 15, enthalten. Zudem kann eine Erkennung von DTMF-Eingaben 16 vorgesehen sein.

Die beschriebene Zuordnung der Wortartkategorie zu den erkannten Worten kann bereits während des Worterkennungsprozesses erfolgen.

Im nächsten Schritt (Schritt 17) werden die erkannten Wörter anhand ihrer Wortkategorien einer VerbalPhrase, d.h. einer auf einem Verb basierenden Phrase, und einer NominalPhrase, d.h. einer auf einem Nomen basierenden Phrase, zugeordnet, vgl. Figur 6.

Danach werden die NomialPhrasen und VerbalPhrasen nach phrasenstrukturellen Gesichtspunkten in Objekten zusammengeführt.

In Schritt 18 werden die Objekte für das Multitasking mit der entsprechenden sprachgesteuerten Anwendung verknüpft.

Jedes Objekt 19 umfasst einen in der Grammar 8 hinterlegten Zielsatz, genauer gesagt ein Satzmodell. Aus Figur 4 geht hervor, dass ein solches Satzmodell z.B. durch eine Wortreihenfolge „Subjekt, Verb, Objekt“ oder „Objekt, Verb, Subjekt“ definiert sein kann. Viele andere Satzbaustrukturen sind in dieser allgemeinen Form in der Grammar 8 hinterlegt. Entsprechen die Wortkategorien der erkannten Wörter 7 der Reihenfolge eines der vorgegebenen Satzmodelle, so werden sie dem zugehörigen Objekt zugeordnet. Der Satz gilt als erkannt. Anders ausgedrückt umfasst jedes Satzmodell eine Anzahl von den verschiedenen Wortkategorien zugeordneten Variablen, die mit den entsprechenden Wortkategorien der erkannten Wörter 7 gefüllt werden.

Das Verfahren bedient sich der traditionellen Grammar Specification Language (GSL), strukturiert die hinterlegten Sätze jedoch in innovativer Weise. Dabei orientiert sie sich an den Regeln der Phrasenstrukturgrammatik und am Konzept einer Generativen Transformationsgrammatik.

Durch die konsequente Anwendung der dort beschriebenen Tiefenstrukturen eines Satzes, insbesondere der Unterscheidung von Nominalphrasen und Verbalphrasen, steht sie der Satzkonstitution der natürlichen Sprache sehr viel näher als die bislang vorherrschenden intuitiven Grammarkonzepte.

Die GT/PS-Grammar basiert somit auf einer theoretischen Modellbildung, die geeignet ist, die abstrakten Prinzipien natürlichsprachlicher Äußerungen zu ermitteln. Auf dem Gebiet moderner Spracherkennungssysteme eröffnet sie erstmals die Möglichkeit, die Abstraktion von Satzbildungsregeln gleichsam umzukehren und als Vorhersage der Äußerungen von Applikationsnutzern zu konkretisieren. Damit wird ein systematischer Zugriff auf Spracherkennungs-Grammars möglich, die bislang stets auf der intuitiven Akkumulation von Beispielsätzen beruhten.

Ein zentrales Merkmal herkömmlicher und GT/PS-Grammars ist die hierarchische Verschachtelung in sogenannte Subgrammars, die einzelne Worte wie Variablen auf der höchsten Ebene zu einem ganzen Satz zusammensetzen. Die GT/PS-Grammar ist in diesem Punkt sehr viel kleiner und hierarchisch viel übersichtlicher als die bisher bekannten Grammars. Im Unterschied zu herkömmlichen Grammars sind in der neuen Grammar fast ausschließlich „sinnvolle“ Sätze hinterlegt, so dass das Maß an Overgeneration, d.h. hinterlegte Sätze, die im natürlichsprachlichen Sinne falsch sind, sinkt. Dies ist wiederum die Voraussetzung für eine verbesserte Erkennerleistung, da die Applikation nur zwischen wenigen hinterlegten Alternativen wählen muss.

Patentansprüche

1. Verfahren zur natürlichen Spracherkennung auf Basis einer Generative Transformations-/Phrasenstruktur-Grammatik, gekennzeichnet durch die Schritte:
 - Analyse einer gesprochenen Phrase auf darin enthaltene Triphone;
 - Bildung von in der gesprochenen Phrase enthaltenen Wörtern aus den erkannten Triphonen mit Hilfe von Lautwortdatenbasen (Dictionaries); und
 - Syntaktische Rekonstruktion der gesprochenen Phrase aus den erkannten Wörtern unter Verwendung eines grammatischen Regelwerks (Grammar).
2. Verfahren nach Anspruch 1, dadurch gekennzeichnet, dass die syntaktische Rekonstruktion der gesprochenen Phrase die Schritte umfasst:
 - Zuordnung der erkannten Wörter zu Wortartenkategorien (Verb, Nomen etc.)
 - Zuordnung der Wortartenkategorien zu Nominalphrasen und Verbalphrasen;
 - Zusammenführung der Nominalphrasen und Verbalphrasen nach syntaktischen Regeln in Objekten unter Vorgabe verschiedene Satzmodelle, wobei die erkannten Wortfolgen mit den vorgegebenen Satzmodellen verglichen werden, wobei im Fall einer Übereinstimmung der Satz als erkannt gilt und eine Aktion in einer sprachgesteuerten Applikation auslöst.
3. Verfahren nach einem der Ansprüche 1 oder 2, dadurch gekennzeichnet, dass jedes Satzmodell eine Anzahl von Wortkategorien zugeordneten Variablen aufweisen, die mit den entsprechenden Wortkategorien der erkannten Wörter gefüllt werden.
4. Verfahren nach einem der Ansprüche 1 bis 3, dadurch gekennzeichnet, dass die zu erkennenden Worte in verschiedene Wortkategorien untergliedert in den Wortdatenbasen vorgehalten werden.

5. Verfahren nach eine der Ansprüche 1 bis 4, dadurch gekennzeichnet, dass die Objekte oder Teile davon mit entsprechenden Aktionsparametern einer sprachgesteuerten Applikation verknüpft werden.

Zusammenfassung

Die Erfindung betrifft ein Verfahren zur natürlichen Spracherkennung auf Basis einer Generativen Transformations-/Phrasenstruktur-Grammatik Die GT-/PS-Grammar. Erfindungsgemäß erfolgt eine Analyse einer gesprochenen Phrase auf darin enthaltene Triphone, eine Bildung von in der gesprochenen Phrase enthaltenen Wörtern aus den erkannten Triphonen mit Hilfe von Lautwortdatenbasen (Dictionaries) und eine syntaktische Rekonstruktion der gesprochenen Phrase aus den erkannten Wörtern unter Verwendung eines grammatischen Regelwerks (Grammar).

Die GT-/PS-Grammar ist ein neuartiges Verfahren, Zielsätze in der Grammar zu hinterlegen. Sie bedient sich der traditionellen Grammar Specification Language (GSL), strukturiert die hinterlegten Sätze jedoch in innovativer Weise. Dabei orientiert sie sich an den Regeln der Phrasenstrukturgrammatik und an Noam Chomskys Konzept einer Generativen Transformationsgrammatik.

1. Triphonanalyse

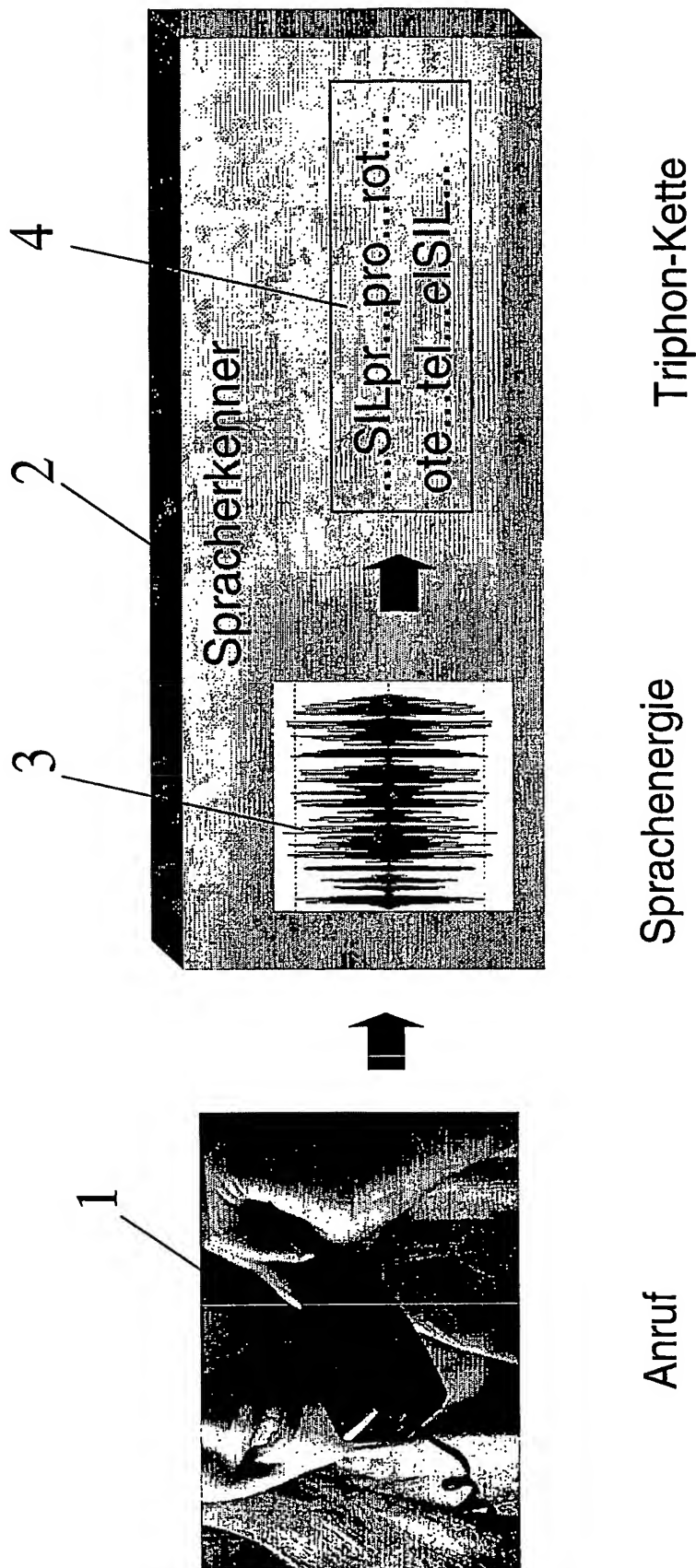


Fig. 1

2. Worterkennung (Dictionaries)

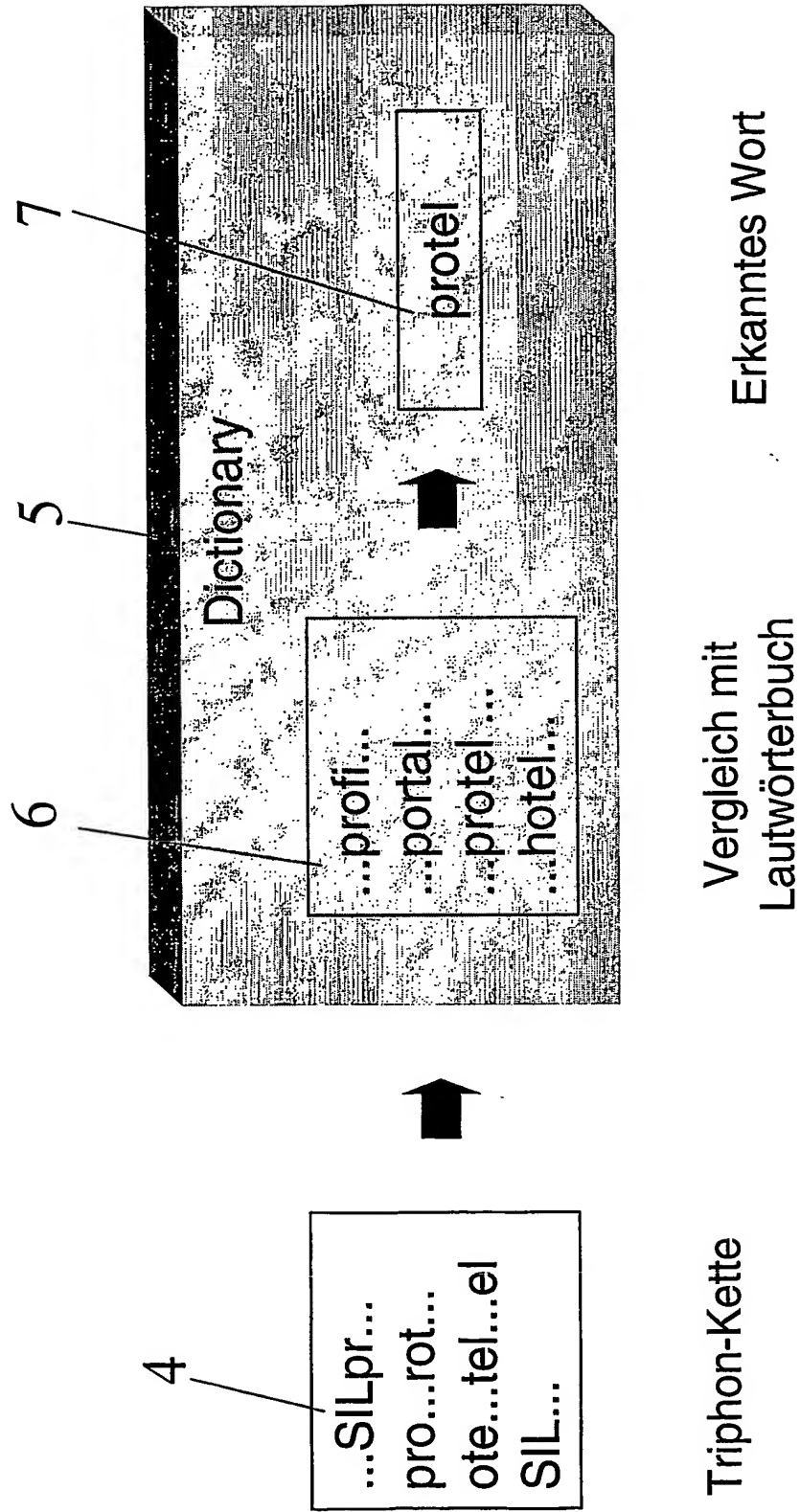


Fig. 2

3. Syntaktische Rekonstruktion (Grammar)

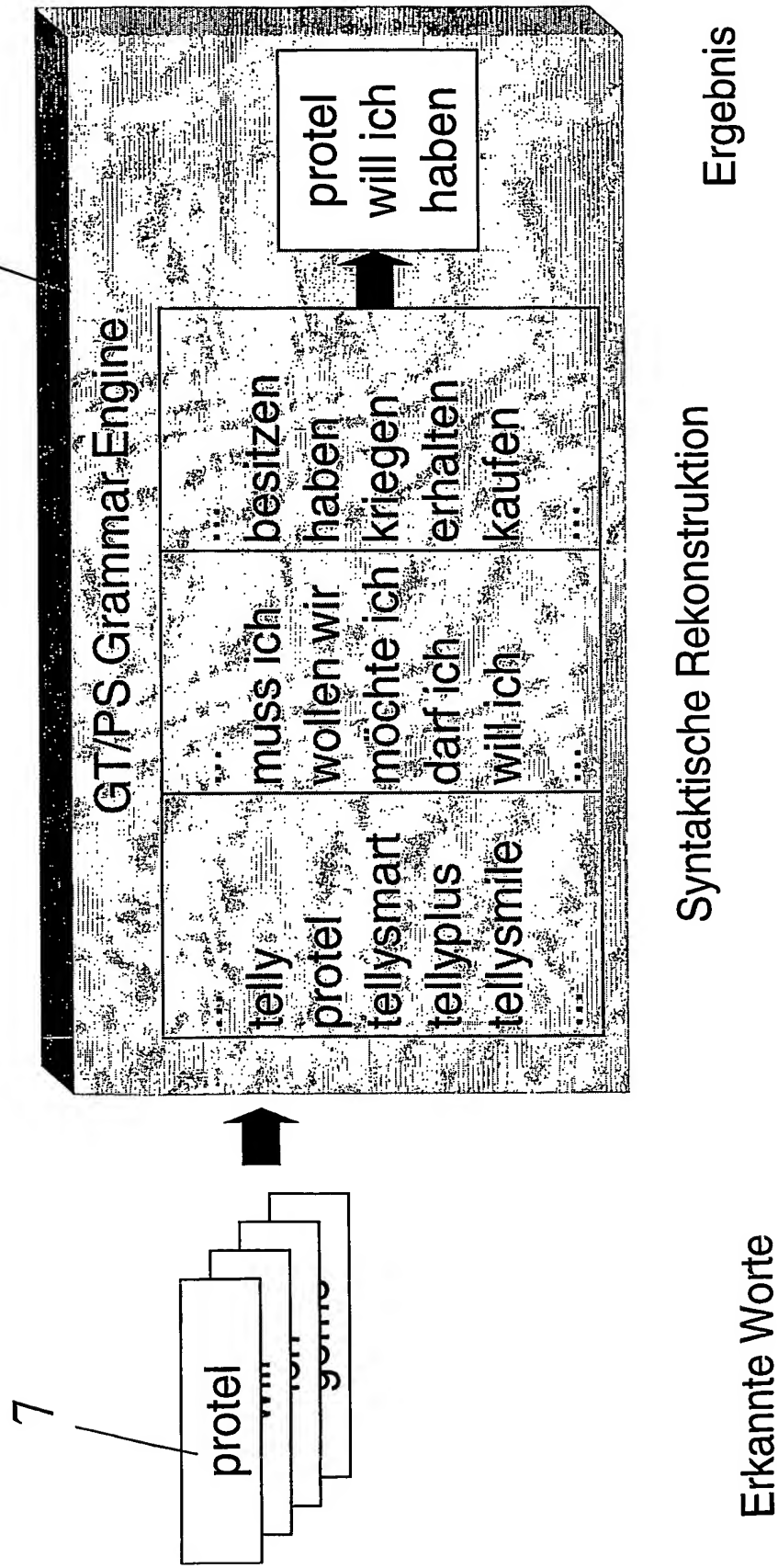


Fig. 3



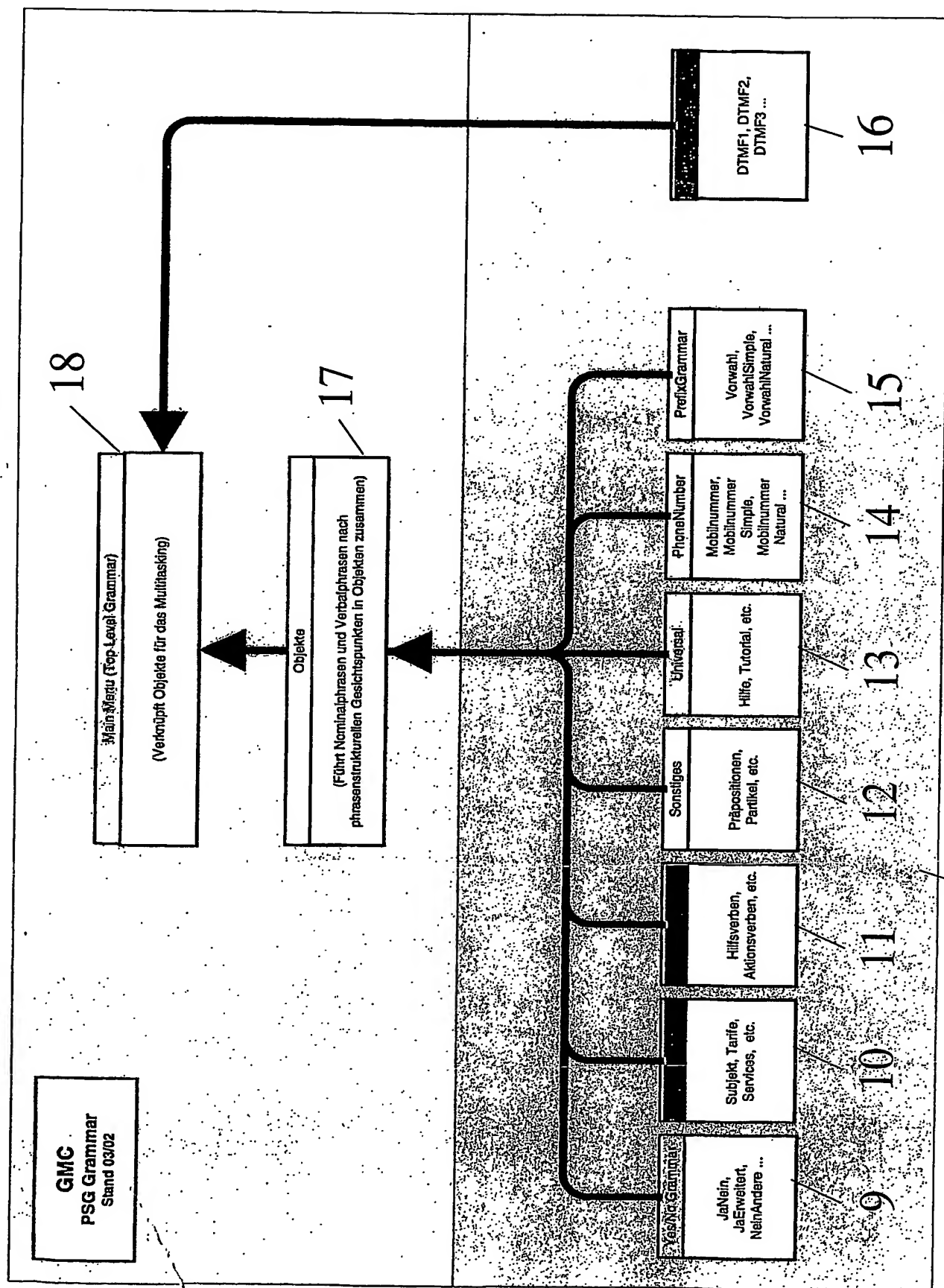


Fig. 5

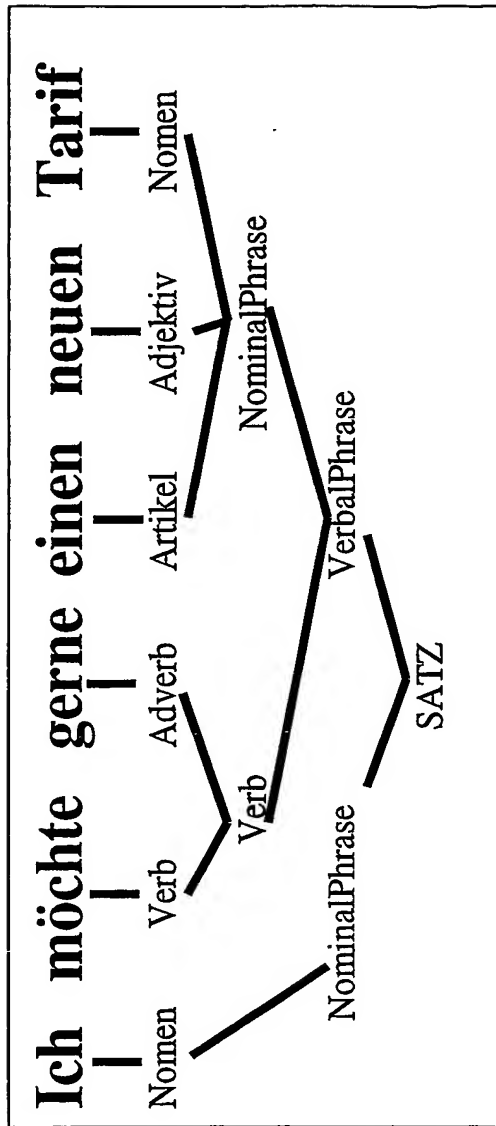


Fig. 6

(Stand der Technik)

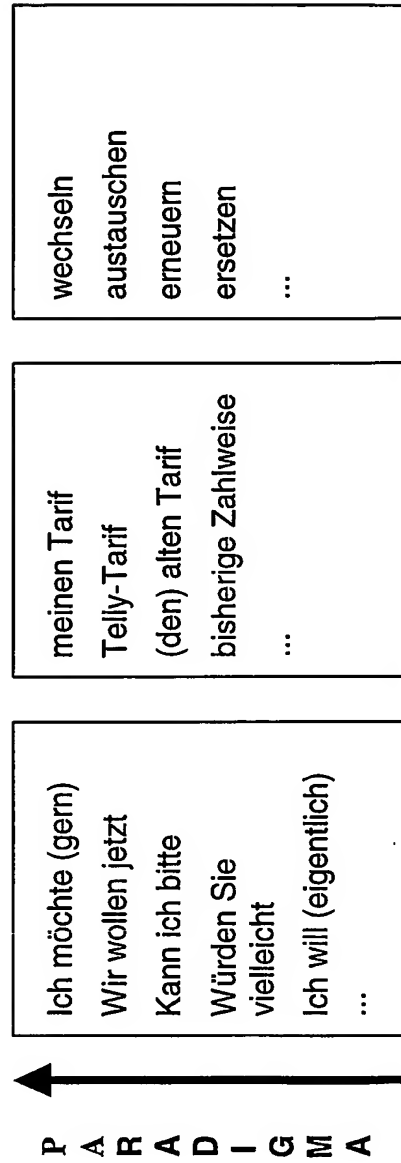


Fig. 7